

Bibliographic note: the material starting w/ the stopping conditions for single-link clustering all comes from Kleinberg, NIPS 2003. (see webpage for link: full bib. info).

3/12/15  
Lecture 13!  
~~Clustering impossibility~~  
~~An impossibility result~~  
~~for clustering~~

Announcement re: grades for AI.

In our explorations of clustering, we've ~~been~~ taking an "optimization function"-based perspective:

- ~~each clustering~~ ~~only~~ ~~could~~

So far: ~~each clustering~~ ~~also~~ ~~could~~ ~~be~~ ~~thought~~ ~~of~~ ~~as~~ ~~an~~ ~~optimization~~ ~~problem~~, and ~~produced~~ ~~produces~~ ~~K~~ ~~clusters~~

- ~~each~~ ~~problem~~ ~~is~~

3 clustering ideas

3 ideas for producing ~~clusters~~ clustering w/ associated optimization function.

~~say, to prevent students feeling like~~

~~say, don't write the rest so that we don't spend too much time in review.~~

- k-means: minimize ~~square~~ within-cluster squared distances.  
can solve approximately
- single-link: maximize the smallest spacing between clusters.  
can solve exactly.
- spectral clustering: find embedding that gives smallest normalized cut.  
~~some~~ approximate solution. >

~~Now, if we have too~~

~~Now, I've only~~

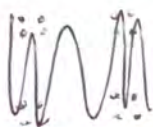
This lecture, going to ~~take a~~ ~~little~~ look @ these premises more carefully:

k fixed so far, but  
shouldn't ~~this~~ ~~k~~ ~~vary~~?  
it ~~depends~~ ~~on~~ ~~the~~ ~~data~~?  
it ~~varies~~ ~~depending~~ ~~on~~ ~~the~~ ~~data~~?

which one best when?

picking

ex: ~~spx we choose~~ ~~k=3~~, and ~~our~~ ~~data~~ ~~looks~~ ~~like~~ ~~this~~:



vs.

"ought" to be 3.

⇒ best to allow # of clusters to vary. "ought" to be 4 clusters

Consider the following single fix to allow ~~adaptive~~ adjusting the # of clusters on the fly:

for given opt. fn, ~~try all possible # of clusters~~ and pick the best clustering  
 get the ~~(approximate)~~ best<sup>↑</sup> for each  $k$ , and then pick the best among those  
 let  $k^* = \#$  of clusters in that idea

- best for ~~min cluster~~ <sup>sum of squared distances</sup>: ~~every point in its own cluster~~ @  $k^* = n$

- best for ~~single-link~~: ~~every point in same cluster~~.

→ ~~maximizing the within-cluster spacing~~  
 → maximizing the between-cluster spacing: farthest 2 points in diff clusters (doesn't matter about anyone else)

(what does this mean if there aren't 2 diff clusters?)

ill-defined for 1 cluster only. But o.w., it works and you just pick the best.

- best for normalized cut:

$$\frac{(\# \text{ of edges cut})}{\sum_j (\# \text{ edges inside } C_j)} : \text{one big cluster. } k^* = 1.$$

and that will be the best two clusters

$k^* = 2$

- maybe best to skip all this: it would require us to write down all the objective functions anyway, and the only "problem" that results from single-link clustering is the ~~problem~~ being ill-defined @  $k=1$  (w.r.t. objective fn).

- could ask as matching ~~cluster~~ clicker.g?

Let  $k^*$  be the # of clusters in the optimal cluster

Let  $k^*$  be the <sup>best</sup> ~~optimal~~  $k$

Let  $k^* = \#$  of clusters in the ~~optimal~~ clustering optimizing the specified function.  
 What is  $k^*$ ? (Assume  $n \gg 2$ )

Q1: for k-means' opt. fn?  
 - if ~~not~~ has

(A)  $k^* = 1$

(B)  $k^* = 2$

(C)  $k^* = 3$

(D)  $k^* = n - 1$

(E)  $k^* = n$

(F) I don't know.

Q2: for normalized cut?

Q3: for single-link? (we ~~do not allow~~  $k^* = 1$ )  
 ignore

Need extra criteria for sensible output.

So quickly! but spend for this number the fix.



So, what else might we try? Need sth about the data points to guide us.  
 to develop intuition, let's ~~take~~ take single-link as a concrete example, and  
 consider how we might get it to stop

clicker q. 4: (diff. possible stopping conditions) for single-link clustering is the "best" way to allow it to ~~create diff~~ produce different #s of clusters?

(A) Instead of stopping cluster merging when there are  $k$  connected components, stop when there are no more diff-cluster pairs <sup>of pairs</sup>  $x_i, x_j$  such that  $d(x_i, x_j) \leq \theta$ . Let  $\theta$  be a fixed positive #, and

(A) ... for fixed threshold  $\theta$ , stop when there are no more diff-cluster pairs with distance  $\leq \theta$ .

(A) ... for fixed threshold  $\theta$ , stop when there are no more diff-cluster pairs with distance  $\geq \theta$ .

(B) ... for fixed  $\theta$  multiplier, stop when there are  $\theta$  clusters that have distance  $\leq \theta \cdot \max_{i,j} d(x_i, x_j)$ .

(C) both are good ~~if you have the inequality in (B) going the wrong way~~. if you changed the " $\geq$ " in (B) to " $\leq$ ".

(D) both would be good if you changed

(C) Exactly one of the ~~" $\leq$ "~~ "Neither are good, but <sup>exactly one is</sup> one would be if you changed a " $\leq$ " to a " $\geq$ " but both would be if you changed

(D) Neither are good, but both would be if you changed both " $\leq$ " to " $\geq$ ".

(E) what's single-link clustering?

(note: this allows extension to  $K \neq 1$ )

~~note that we lose our guarantee~~

Other examples:

For other

~~could try to do sth similar~~

skip

- Other alternatives:

- add penalty term to avoid overfitting (K-means)

~~nearest neighbor~~

- eigen-gap in the spectrum (list of eigenvalues)

~~All these~~ NOW

All these questions ~~are motivated~~ about trying to know how many clusters ~~we~~ one wants. All these questions were motivated by the principle that one really can't a priori know how many clusters ~~we~~ one wants.

We might ~~encode this~~ enshrine this principle by giving it a name:  
 \* Formalizing ~~the~~ the idea that ~~# of clusters need not be fixed~~ any # of

Let  $X = \{x_1, \dots, x_n\}$  our points.

$d: X \times X \rightarrow \mathbb{R}$  a distance function.

And this principle operated "above" the level of specific optimization fns.

Idea: ~~are there principle~~

Principle: whatever optimization function or alg we choose, it should be able to produce any possible clustering:

We can make this principle even more general, and say that not only should the number of clusters be allowed to vary over all possible #'s, but that the clusterings themselves should be allowed to vary over all possible partitions.  
 Our work suggests:  
 A meta-principle (rather than an opt. fn):

~~Ex: set~~

for convenience let's number points, and say  $X = \{1, 2, \dots, n\}$ .  
 For any partition  $\Pi$  of  $X$ , we can set up <sup>the</sup> "natural" clustering points when that's "natural" clustering.

ex:  $\Pi = \{1, 2, 3\}, \{4, 5, 6\}$

ex:  $\Pi = \{1, 2, 3\}, \{3, 4, 5\}, \{6\}$



distances

$d: X \times X \rightarrow \mathbb{R}^{\geq 0}$

So, define a clustering fn as:  $f$ : input = ~~distance fn~~ distance fn  $d(x, y)$  on points  
 output: a partition of  $X$ .

Richness property:  $\text{Range}(f) = \text{set of all partitions of } X$ .

<whatever an alg is, it can't rule out some clustering a priori >  
 <note that all our fixed-k algs violated this.>

\* [Cornell's] Jon

\* Note: (Kleinberg's insight: meta-analysis: level above specific opt. fns, Note: all our fixed-k algs violated this. ~~But~~ we talked about fixes. <sthg we'll get back to later... >



Now, while we're @ it, let's look @ other reasonable properties.

If the only input we have is the distance function onto the non-negative #'s - note there are no "units", then it shouldn't matter what the units are.

Scale-invariance property: (no "units")

for any distance  $f$  and any  $\alpha > 0$ , let  $d_\alpha(t,s) = \alpha d(t,s)$  ("scale it").

$f(d)$  =  $f(\frac{d}{\alpha})$   
~~partitioned you get result partition.~~  
~~resulting partition is the distance~~  
 <If you spread all the distances by the same amount, get same result>  
 <all the distances have the same relative relationships to them>  
 <seems like k-means, single-link, (tie!)>

<drawing a picture actually makes it seem more suspect>

Consistency property: [if you shrink "w/in-cluster distances"; expand "between-cluster distances",  $f$  gives the same partition out.]

<data is only becoming even more well-clustered>

- this "defn" requires introducing a specific partition:

let  $d$  be a distance function:

$P = f(d)$  the output partition

and  $d' = \alpha d$  distance function where  $t,s$  are ~~in~~ in the same cluster in  $P$   
 $\Rightarrow d'(t,s) \leq d(t,s)$  (~~in~~  $(t,s) \in P$ )

"shrinks inter-cluster distances"

and  $t,s$  in diff clusters in  $P$  ( $(t,s) \notin P$ )

$\Rightarrow d'(t,s) \geq d(t,s)$

"expands b/w cluster distances"

~~f satisfies~~

Then,  $f(d') = P (= f(d))$

P: if look too similar

Q5: which does our fixed-k single-link clustering satisfy?

(D) = B & C

Q

Q6: which does our distance- $\theta$  single-link clustering satisfy?

(D) = A & C

Q7: which does ~~our~~ our fixed-k-means alg satisfy?

- maybe requires pf! (B) & (C).

<not (A) is obvious>

- (A) <sup>just</sup> richness
- (B) <sup>just</sup> scale-invariance
- (C) <sup>just</sup> consistency
- (D) exactly 2
- (E) all 3

OK: so, is there an alg that satisfies all 3 properties?

Well, clearly spectral clustering, right? o.w., why would we have ~~done~~ presented it last?

outline:

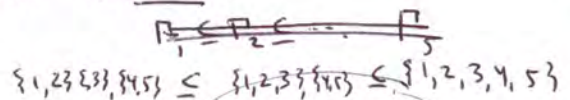
thm 3.1: spsc clustering fn  $f$  satisfies scale-invariance & consistency.

Then,  $\text{Range}(f)$  does not contain any pair of partitions where one refines the other.

~~We say that Range~~

That is,  $\text{Range}(f)$  is an anti-chain. [vocab item.]

not only is it not a chain:



there isn't even any links (a pair)

thm 3.2: for every antichain  $A$ , there is a clustering fn  $f$  s.t.  $\text{Range}(f) = A$  and  $f$  satisfies scale-invariance & consistency.

~~Since the set of all partitions is not an anti-chain~~

~~Since~~

But, the set of all partitions is not an antichain.

$\Rightarrow$  can't have both richness & the other two simultaneously.

(so, no, spectral clustering doesn't satisfy all three)

~~which explain~~

an expression of fundamental tradeoffs in clustering.

formal

<so instead of fighting over which opt. fn is best realize each is giving up something>



let's do 3.2:

shows how we can satisfy scalability; consistency ... as long as we restrict the possible outputs.  
uses our k-means opt-fn  
shows

~~Let  $\mathcal{P}$  in an anti-chain of partitions.~~

let  $f$  output the partition  $P \in \mathcal{A}$  that minimizes:

$$\Phi_d(P) = \sum_{\substack{t,s \text{ in} \\ \text{same cluster in } P}} d(t,s)$$

(recall: w/in cluster scatter, for  $d =$  Euc. dist. squared)

• b/c  $\mathcal{A}$  an anti-chain, all-singleton clusters ~~is allowed~~  $\notin \mathcal{A}$ .

scale-invariance?  $\checkmark$  spreading all the distances by same amt ~~means~~ would still yield same minimum.

~~$\Phi_d(P)$  has same min as~~ would just have the scaling factor here.

$\text{Range}(f) \stackrel{c}{=} \mathcal{A}$ ?  $\checkmark$  by ~~constructive~~ defn.

prove:  $\mathcal{A} \subseteq \text{Range}(f)$ . to get that we can get this "pseudo-richness":  
i.e.: show: for any  $\hat{P} \in \mathcal{A}$ , there's a distance matrix  $d$  so that  $f(d) = \hat{P}$ .

just make all in-cluster distances really small, ex:  $d(t,s) < \frac{1}{n^3}$  if  $t,s \in P$ .  
 ... ~~blow~~ ... relatively big, ex:  $d(t,s) > 1$ .

so,  $\Phi_d(P) = \sim n^2$  sum of  $< \frac{1}{n^3}$  terms,  $< 1$ .

$\Rightarrow$  for  $f$  to ~~pick~~ <sup>output</sup> something other than  $P$ , that "something" must also  $\Phi_d(P')$  must be  $< 1$ .

$\Rightarrow$  ~~can~~  $P'$  must refine  $P$ . (o.w, would have an edge)

clicker Q8:  
for blow-cluster distans, set  $d(t,s) =$   
(A)  $\frac{1}{n}$  (B)  $\frac{1}{n^2}$  (C) 1  
(D) I DK.

after stuff below  $\rightarrow$

• need  $\hat{P}$  to be ~~the~~ <sup>to be argmin</sup> ~~minimizer~~ of  $\Phi_d(P)$   $P \in \mathcal{A}$ . (Should use the anti-chain fact)

• if can show any  $P'$  w/  $\Phi_d(P') < \Phi_d(P)$ ,  $P'$  is a refinement of  $\hat{P}$ , done. (b/c then  $\hat{P} \notin \mathcal{A}$ ,

\* idea: force any 'minimizer' to not have any blow-cluster  $\hat{P}$  pairs  $t,s$  w/in a clstr.

$\Rightarrow$

consistency? (let's skip)

to show: for a ~~some~~ <sup>given</sup>  $d$ , and so let  $\hat{P} = f(d)$ .  
 take  $d'$  s.t. ~~in~~ <sup>any</sup>  $w$  in cluster for  $P$  squished,  $b$  in cluster for  $P$  increased.

Show:  $f(d) = f(d')$ ,  $\hat{P} = \hat{P}$

~~we know  $\hat{P} = \arg \min_P \phi_d(P)$ , show  $\hat{P} = \arg \min_P \phi_{d'}(P)$~~

⊙ for any other  $P'$ ,  $\phi_d(P') \geq \phi_d(\hat{P}) \geq \phi_{d'}(\hat{P})$   
 at least  $\phi_d(\hat{P})$  smallest for  $d$ .

We know  $\phi_d(P') \geq \phi_d(\hat{P})$  b/c of defn of  $f$ .

want:  $\phi_{d'}(P') \geq \phi_{d'}(\hat{P})$ .

so take  $\phi_d(P') - \phi_{d'}(P')$  vs.  $\phi_d(\hat{P}) - \phi_{d'}(\hat{P})$ .

$$= \sum_{\substack{\text{w/in cluster} \\ \text{pairs for } P'}} [d(s,t) - d'(s,t)]$$

$$\leq \sum_{\substack{s,t \in P' \\ \text{and} \\ s,t \in \hat{P}}} [d(s,t) - d'(s,t)]$$

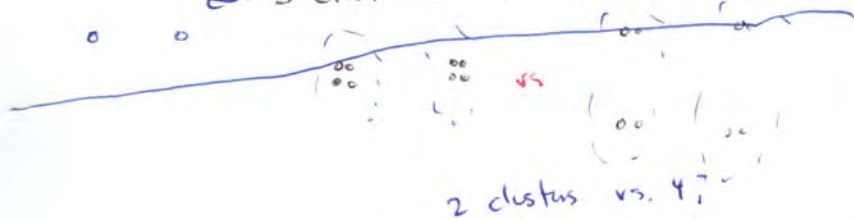
added this here:  
 for edges ~~not in~~  $\hat{P}$  but in  $P$  within,  $d'(s,t) > d(s,t)$  b/c of  $d'$  being a  $P$ -transform

$$\leq \sum_{s,t \in \hat{P}} [d(s,t) - d'(s,t)]$$

added more edges that are  $\sim$  in  $\hat{P}$ , so non-neg. quant is added.

$$= \phi_d(\hat{P}) - \phi_{d'}(\hat{P}), \text{ as desired.}$$

⊙ 3 clusters or one? Tradeoff



The point is not that the properties are "bulletproof".  
 The point is to be able to formalize the tradeoffs.